Carnegie Mellon University ADVANCED DATABASE SYSTEMS Parallel Join Algorithms (Hashing)

 $\overline{}$

#

ctu

()

#FREF PAVLO

@Andy_Pavlo // 15-721 // Spring 2019

TODAY'S AGENDA

Background Parallel Hash Join Hash Functions Hashing Schemes Evaluation





PARALLEL JOIN ALGORITHMS

Perform a join between two relations on multiple threads simultaneously to speed up operation.

- Two main approaches:
- \rightarrow Hash Join
- \rightarrow Sort-Merge Join

We won't discuss nested-loop joins...



OBSERVATION

Many OLTP DBMSs don't implement hash join.

But an **index nested-loop join** with a small number of target tuples is more or less equivalent to a hash join.



HASHING VS. SORTING

- **1970s** Sorting
- **1980s** Hashing
- 1990s Equivalent
- 2000s Hashing
- 2010s Hashing (Partitioned vs. Non-Partitioned)
 2020s ???



PARALLEL JOIN ALGORITHMS



JOIN ALGORITHM DESIGN GOALS

Goal #1: Minimize Synchronization → Avoid taking latches during execution.

Goal #2: Minimize CPU Cache Misses

 \rightarrow Ensure that data is always local to worker thread.



IMPROVING CACHE BEHAVIOR

Factors that affect cache misses in a DBMS:

- \rightarrow Cache + TLB capacity.
- \rightarrow Locality (temporal and spatial).

Non-Random Access (Scan):

- \rightarrow Clustering to a cache line.
- \rightarrow Execute more operations per cache line.

Random Access (Lookups):

 \rightarrow Partition data to fit in cache + TLB.

Source: Johannes Gehrke



PARALLEL HASH JOINS

Hash join is the most important operator in a DBMS for OLAP workloads.

It's important that we speed it up by taking advantage of multiple cores.

 \rightarrow We want to keep all of the cores busy, without becoming memory bound



HASH JOIN (R⊠S)

Phase #1: Partition (optional)

 \rightarrow Divide the tuples of **R** and **S** into sets using a hash on the join key.

Phase #2: Build

 \rightarrow Scan relation **R** and create a hash table on join key.

Phase #3: Probe

 \rightarrow For each tuple in **S**, look up its join key in hash table for **R**. If a match is found, output combined tuple.



PARTITION PHASE

Split the input relations into partitioned buffers by hashing the tuples' join key(s).

- \rightarrow Ideally the cost of partitioning is less than the cost of cache misses during build phase.
- → Sometimes called *hybrid hash join*.

Contents of buffers depends on storage model:

- \rightarrow **NSM**: Either the entire tuple or a subset of attributes.
- \rightarrow **DSM**: Only the columns needed for the join + offset.



PARTITION PHASE

Approach #1: Non-Blocking Partitioning

- \rightarrow Only scan the input relation once.
- \rightarrow Produce output incrementally.

Approach #2: Blocking Partitioning (Radix)

- \rightarrow Scan the input relation multiple times.
- \rightarrow Only materialize results all at once.
- \rightarrow Sometimes called *radix hash join*.



NON-BLOCKING PARTITIONING

Scan the input relation only once and generate the output on-the-fly.

Approach #1: Shared Partitions

- \rightarrow Single global set of partitions that all threads update.
- \rightarrow Have to use a latch to synchronize threads.

Approach #2: Private Partitions

- \rightarrow Each thread has its own set of partitions.
- \rightarrow Have to consolidate them after all threads finish.



Data Table



Data Table



Partitions



CARNEGIE MELLON DATABASE GROUP













CARNEGIE MELLON DATABASE GROUP



CARNEGIE MELLON DATABASE GROUP



RADIX PARTITIONING

Scan the input relation multiple times to generate the partitions.

Multi-step pass over the relation:

- \rightarrow Step #1: Scan R and compute a histogram of the # of tuples per hash key for the <u>radix</u> at some offset.
- → Step #2: Use this histogram to determine output offsets by computing the prefix sum.
- \rightarrow Step #3: Scan R again and partition them according to the hash key.

RADIX

The radix is the value of an integer at a particular position (using its base).





RADIX

The radix is the value of an integer at a particular position (using its base).





RADIX

The radix is the value of an integer at a particular position (using its base).















CARNEGIE MELLON DATABASE GROUP

RADIX PARTITIONS

Step #1: Inspect input, create histograms







CARNEGIE MELLON DATABASE GROUP



Step #2: Compute output offsets





CARNEGIE MELLON DATABASE GROUP
Step #3: Read input and partition Partition 0, CPU 0 **[#**, 07 пп **(#**_p) 8 07 **(#**_p) Partition 0, CPU 1 03 Ο hash_p(key) Partition 0: 2 [**#**_p Partition 1, CPU 0 0 18 Partition 1: 2 [**#**_p] 19 03 пп [**#**_p Partition 1, CPU 1 **(#**_p) 15 5 Partition 0: 1 [**#**_p 10 Partition 1: 3

CARNEGIE MELLON DATABASE GROUP

Recursively repeat until target number of partitions have been created



Recursively repeat until target number of partitions have been created



Recursively repeat until target number of partitions have been created



BUILD PHASE

The threads are then to scan either the tuples (or partitions) of R.

For each tuple, hash the join key attribute for that tuple and add it to the appropriate bucket in the hash table.

 \rightarrow The buckets should only be a few cache lines in size.



HASH TABLE

Design Decision #1: Hash Function

- \rightarrow How to map a large key space into a smaller domain.
- \rightarrow Trade-off between being fast vs. collision rate.

Design Decision #2: Hashing Scheme

- \rightarrow How to handle key collisions after hashing.
- \rightarrow Trade-off between allocating a large hash table vs. additional instructions to find/insert keys.



HASH FUNCTIONS

We don't want to use a cryptographic hash function for our join algorithm.

We want something that is fast and will have a low collision rate.

- \rightarrow **Best Speed:** Always return '1'
- → **Best Collision Rate:** Perfect hashing



HASH FUNCTIONS

<u>CRC-32</u> (1975)

 \rightarrow Used in networking for error detection.

MurmurHash (2008)

 \rightarrow Designed to a fast, general purpose hash function.

Google CityHash (2011)

 \rightarrow Designed to be faster for short keys (<64 bytes).

Google FarmHash (2014)

 \rightarrow Newer version of CityHash with better collision rates.

<u>CLHash</u> (2016)

 \rightarrow Fast hashing function based on <u>carry-less multiplication</u>.



HASH FUNCTION BENCHMARK



HASH FUNCTION BENCHMARK



HASHING SCHEMES

Approach #1: Chained Hashing

Approach #2: Linear Probe Hashing

Approach #3: Robin Hood Hashing

Approach #4: Cuckoo Hashing



CHAINED HASHING

Maintain a linked list of <u>buckets</u> for each slot in the hash table.

- Resolve collisions by placing all elements with the same hash key into the same bucket.
- \rightarrow To determine whether an element is present, hash to its bucket and scan for it.
- \rightarrow Insertions and deletions are generalizations of lookups.



CHAINED HASHING





Single giant table of slots.

Resolve collisions by linearly searching for the next free slot in the table.

- → To determine whether an element is present, hash to a location in the index and scan for it.
- \rightarrow Have to store the key in the index to know when to stop scanning.
- \rightarrow Insertions and deletions are generalizations of lookups.















CARNEGIE MELLON DATABASE GROUP

LINEAR PROBE HASHING





30

30



OBSERVATION

To reduce the # of wasteful comparisons during the join, it is important to avoid collisions of hashed keys.

This requires a chained hash table with $\sim 2x$ the number of slots as the # of elements in **R**.



Variant of linear hashing that steals slots from "rich" keys and give them to "poor" keys.

- \rightarrow Each key tracks the number of positions they are from where its optimal position in the table.
- → On insert, a key takes the slot of another key if the first key is farther away from its optimal position than the second key.















CARNEGIE MELLON DATABASE GROUP

ROBIN HOOD HASHING





33









CUCKOO HASHING

Use multiple tables with different hash functions.

- \rightarrow On insert, check every table and pick anyone that has a free slot.
- → If no table has a free slot, evict the element from one of them and then re-hash it find a new location.

Look-ups are always O(1) because only one location per hash table is checked.



CUCKOO HASHING



CUCKOO HASHING








CARNEGIE MELLON DATABASE GROUP





 $\frac{\text{Insert } X}{\text{hash}_{1}(X) \text{ hash}_{2}(X)}$

Insert Y hash₁(Y) hash₂(Y)

Insert Z ∽ hash₁(Z) hash₂(Z) ∽







Hash Table #1



 $\frac{\text{Insert } X}{\text{hash}_1(X) \quad \text{hash}_2(X)}$

Insert Y hash₁(Y) hash₂(Y)

 $\frac{\text{Insert } Z}{\text{hash}_1(Z) \quad \text{hash}_2(Z)} \checkmark$

Hash Table #2









 $\frac{\text{Insert } X}{\text{hash}_1(X) \text{ hash}_2(X)}$

Insert Y hash₁(Y) hash₂(Y)

Insert Z $hash_1(Z)$ $hash_2(Z)$ $-hash_1(Y)$

Hash Table #2





Hash Table #1



 $\frac{\text{Insert } X}{\text{hash}_1(X) \quad \text{hash}_2(X)}$

Insert Y hash₁(Y) hash₂(Y)

 $\begin{array}{c}
 Insert Z \\
 hash_1(Z) & hash_2(Z) \\
 hash_1(Y)
\end{array}$

Hash Table #2





Hash Table #1



 $\begin{array}{c} \text{Insert } X\\ \text{hash}_{1}(X) & \text{hash}_{2}(X) \end{array}$

Insert Y hash₁(Y) hash₂(Y)

Insert Z $hash_1(Z) \quad hash_2(Z)$ $hash_1(Y)$ $hash_2(X)$







Threads have to make sure that they don't get stuck in an infinite loop when moving keys.

If we find a cycle, then we can rebuild the entire hash tables with new hash functions.

- \rightarrow With <u>two</u> hash functions, we (probably) won't need to rebuild the table until it is at about 50% full.
- → With <u>three</u> hash functions, we (probably) won't need to rebuild the table until it is at about 90% full.



PROBE PHASE

For each tuple in **S**, hash its join key and check to see whether there is a match for each tuple in corresponding bucket in the hash table constructed for **R**.

- \rightarrow If inputs were partitioned, then assign each thread a unique partition.
- \rightarrow Otherwise, synchronize their access to the cursor on **S**



MICRO ADAPTIVITY IN VECTORWISE

IGMOD 2013

ABASE GROUP

PROBE PHASE - BLOOM FILTER

Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

- \rightarrow Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.
- → Sometimes called *sideways information passing*.



38

MICRO ADAPTIVITY IN VECTORWISE

IGMOD 2013

TABASE GROUP

PROBE PHASE - BLOOM FILTER

Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

- \rightarrow Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.
- → Sometimes called *sideways information passing*.



PROBE PHASE - BLOOM FILTER

Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

- \rightarrow Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.
- → Sometimes called *sideways information passing*.



Bloom Filter

MICRO ADAPTIVITY IN VECTORWISE

IGMOD 2013

ABASE GROUP

PROBE PHASE - BLOOM FILTER

Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

- \rightarrow Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.
- → Sometimes called *sideways information passing*.



HASH JOIN VARIANTS

	No-P	Shared-P	Private-P	Radix
Partitioning	No	Yes	Yes	Yes
Input scans	0	1	1	2
Sync during partitioning	_	Spinlock per tuple	Barrier, once at end	Barrier, 4 · #passes
Hash table	Shared	Private	Private	Private
Sync during build phase	Yes	No	No	No
Sync during probe phase	No	No	No	No

ATABASE GROUP

BENCHMARKS

Primary key – foreign key join \rightarrow Outer Relation (Build): 16M tuples, 16 bytes each \rightarrow Inner Relation (Probe): 256M tuples, 16 bytes each

Uniform and highly skewed (Zipf; s=1.25)

No output materialization

DESIGN AND EVALUATION OF MAIN MEMORY HASH JOIN ALGORITHMS FOR MULTI-CORE CPUS SIGMOD 2011

HASH JOIN - UNIFORM DATA SET

Intel Xeon CPU X5650 @ 2.66GHz 6 Cores with 2 Threads Per Core



HASH JOIN - SKEWED DATA SET

Intel Xeon CPU X5650 @ 2.66GHz 6 Cores with 2 Threads Per Core



OBSERVATION

We have ignored a lot of important parameters for all of these algorithms so far.

- \rightarrow Whether to use partitioning or not?
- \rightarrow How many partitions to use?
- \rightarrow How many passes to take in partitioning phase?

In a real DBMS, the optimizer will select what it thinks are good values based on what it knows about the data (and maybe hardware).



RADIX HASH JOIN - UNIFORM DATA SET

Intel Xeon CPU X5650 @ 2.66GHz Varying the # of Partitions



RADIX HASH JOIN - UNIFORM DATA SET

Intel Xeon CPU X5650 @ 2.66GHz Varying the # of Partitions



EFFECTS OF HYPER-THREADING



PARTING THOUGHTS

On modern CPUs, a simple hash join algorithm that does not partition inputs is competitive.

There are additional vectorization execution optimizations that are possible in hash joins that we didn't talk about. But these don't really help...



NEXT CLASS

Parallel Sort-Merge Joins

